

K-Mean Class Family

Andres Mendez-Vazquez
Cinvestav Guadalajara
amendez@gdl.cinvestav.mx

Abstract

I am putting together a class in python that implement several of the k-means algorithms. This is still a work under progress so excuses any mistake.

1. Introduction

This class implements a series of classic algorithms from the k-mean family. The algorithms being implemented are:

- K-Means
- K-Medians
- K-Centers
- Fast K-Medoids
- Fuzzy C-Means

2. Auxiliary Functions

2.1. Canopy Initialization

The Centroid Initialization

2.2. KMetric

Allows to select three different metrics:

- Euclidean Metric - The only differentiable metric in the list (a.k.a. No singular points).

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{(\mathbf{x} - \mathbf{y})^T (\mathbf{x} - \mathbf{y})} \quad (1)$$

- Manhattan Metric

$$dist(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i| \quad (2)$$

- Chebyshev Metric

$$dist(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i| \quad (3)$$

2.3. set_k

It allows to reset the number of clusters for different experiments

3. Implemented Algorithms

3.1. K-Means

This is the classic algorithm was first proposed by Stuart Lloyd in 1957 as a technique for pulse-code modulation. Yes! Signal processing.

It implements an algorithm to minimize the following cost function

$$\sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 = \sum_{k=1}^N \sum_{i:\mathbf{x}_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (4)$$

3.1.1. The Main Algorithm

The main algorithm is described below.

K-means(X, k)

1. Randomly choose K data points (seeds) to be the initial centroids, cluster centers,

$$\bullet \{ \mathbf{v}_1, \dots, \mathbf{v}_k \}$$

2. Assign each data point to the closest centroid

$$\bullet c_i = \arg \min_j \{ dist(\mathbf{x}_i - \mathbf{v}_j) \}$$

3. Re-compute the centroids using the current cluster memberships.

$$\bullet \mathbf{v}_j = \frac{\sum_{i=1}^n I(c_i = j) \mathbf{x}_i}{\sum_{i=1}^n I(c_i = j)}$$

4. If a convergence criterion is not met, go to 2.

3.2. K-Medians

The k-Median algorithm is a concave minimization with cost function

$$\min_{C,D} \sum_{i=1}^m \min_{l=1,\dots,k} \|D_{il}\|_p \text{ s.t. } -D_{il} \leq A^T - C_l \leq D_{il}, i = 1, \dots, m, l = 1, \dots, k \quad (5)$$

3.3. K-Centers

The K-center criterion partitions the points into k clusters so as to minimize the maximum distance of any point to its cluster center. It minimizes the worst case distance to the centroids

$$\min_S \max_{k=1,\dots,K} \max_{i:x_i \in C_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)^T (\mathbf{x}_i - \boldsymbol{\mu}_k) \quad (6)$$

K-centers Algorithm

- Step 1
 - Randomly select an object \mathbf{x}_j from S , let $\mathbf{h}_1 = \mathbf{x}_j$, $H = \{\mathbf{h}_1\}$.
 - It does not matter how \mathbf{x}_j is selected.
- Step 2
 - For $j = 1$ to n :
 - * $dist(\mathbf{x}_i) = L(\mathbf{x}_i, \mathbf{h}_1)$.
 - * $cluster(\mathbf{x}_i)$
- Step 3
 - For $i = 2$ to k
 1. $D = \max_{\mathbf{x}_j: \mathbf{x}_j \in S-H} dist(\mathbf{x}_j)$
 2. Choose $\mathbf{h}_i \in S - H$ such that $dist(\mathbf{h}_i) == D$
 3. $H = H \cup \{\mathbf{h}_i\}$
 4. for $j = 1$ to N
 - * if $L(\mathbf{x}_j, \mathbf{h}_i) \leq dist(\mathbf{x}_j)$
 - $dist(\mathbf{x}_j) = L(\mathbf{x}_j, \mathbf{h}_i)$
 - $cluster(\mathbf{x}_j) = i$

3.4. K-Medoids

Similar to K-Means, but instead of using the means as centroid some k elements of the data sets. Then, new k centers are chosen to see if they minimize the total distances to the elements:

$$\sum_{i=1}^N \sum_{j=1}^d |x_{ij} - c_j^i| \quad (7)$$

3.5. Fuzzy C-Means

Here, a relaxation of the K-Means allows to obtain a new cost function

$$J_m(\mathcal{S}) = \sum_{k=1}^N \sum_{i=1}^C [A_i(\mathbf{x}_k)]^m \|\mathbf{x}_k - \mathbf{v}_i\|^2 \quad (8)$$

Under the constraints:

- $A_i(\mathbf{x}_k) \in [0, 1], \text{ for } 1 \leq k \leq N \text{ and } 1 \leq i \leq C$
- $\sum_{i=1}^C A_i(\mathbf{x}_k) = 1, \text{ for } 1 \leq k \leq N$
- $0 < \sum_{k=1}^N A_i(\mathbf{x}_k) < n, \text{ for } 1 \leq i \leq C$
- $m > 1$

Fuzzy-C-Mean Algorithm

1. Let $t = 0$.
2. Select an initial fuzzy pseudo-partition.
3. Calculate the initial C cluster centers using,

$$v_i^{(t)} = \frac{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m \mathbf{x}_k}{\sum_{k=1}^N A_i^{(t)}(\mathbf{x}_k)^m}. \quad (9)$$

4. Update for each x_k the membership function by

- Case I: $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 > 0$ for all $i \in \{1, 2, \dots, C\}$ then

$$A_i^{(t+1)}(\mathbf{x}_k) = \frac{1}{\left[\sum_{j=1}^C \left\{ \frac{\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2}{\|\mathbf{x}_k - \mathbf{v}_j^{(t)}\|^2} \right\}^{\frac{1}{m-1}} \right]} \quad (10)$$

- Case II: $\|\mathbf{x}_k - \mathbf{v}_i^{(t)}\|^2 = 0$ for some $i \in I \subseteq \{1, 2, \dots, C\}$ then define $A_i^{(t+1)}(\mathbf{x}_k)$ by any non-negative number such that $\sum_{i \in I} A_i^{(t+1)}(\mathbf{x}_k) = 1$ and $A_i^{(t+1)}(\mathbf{x}_k) = 0$ for $i \notin I$.

5. If $|\mathcal{S}^{(t+1)} - \mathcal{S}^{(t)}| = \max_{i,k} |A_i^{(t+1)}(\mathbf{x}_k) - A_i^{(t)}(\mathbf{x}_k)| \leq \epsilon$ stop; otherwise increase t and go to step 2.